



**University of  
Zurich**<sup>UZH</sup>

**Zurich Open Repository and  
Archive**

University of Zurich  
University Library  
Strickhofstrasse 39  
CH-8057 Zurich  
[www.zora.uzh.ch](http://www.zora.uzh.ch)

---

Year: 2020

---

## **Recommending investors for new startups by integrating network diffusion and investors' domain preference**

Xu, Shuqi ; Zhang, Qianming ; Lü, Linyuan ; Mariani, Manuel

**Abstract:** Over the past decade, many startups have sprung up, which create a huge demand for financial support from venture investors. However, due to the information asymmetry between investors and companies, the financing process is usually challenging and time-consuming, especially for the startups that have not yet obtained any investment. Because of this, effective data-driven techniques to automatically match startups with potentially relevant investors would be highly desirable. Here, we analyze 34,469 valid investment events collected from [www.itjuzi.com](http://www.itjuzi.com) and consider the cold-start problem of recommending investors for new startups. We address this problem by constructing different tripartite network representations of the data where nodes represent investors, companies, and companies' domains. First, we find that investors have strong domain preferences when investing, which motivates us to introduce virtual links between investors and investment domains in the tripartite network construction. Our analysis of the recommendation performance of diffusion-based algorithms applied to various network representations indicates that prospective investors for new startups are effectively revealed by integrating network diffusion processes with investors' domain preference.

DOI: <https://doi.org/10.1016/j.ins.2019.11.045>

Posted at the Zurich Open Repository and Archive, University of Zurich

ZORA URL: <https://doi.org/10.5167/uzh-178484>

Journal Article

Accepted Version



The following work is licensed under a Creative Commons: Attribution-NonCommercial-NoDerivatives 4.0 International (CC BY-NC-ND 4.0) License.

Originally published at:

Xu, Shuqi; Zhang, Qianming; Lü, Linyuan; Mariani, Manuel (2020). Recommending investors for new startups by integrating network diffusion and investors' domain preference. *Information Sciences*, 515:103-115.

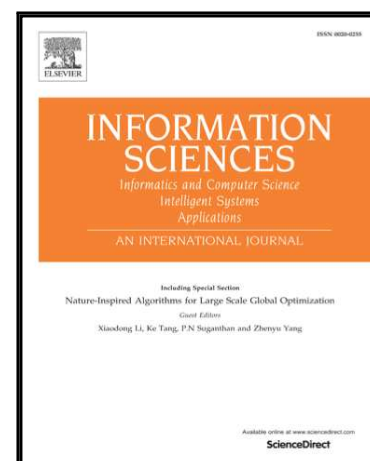
DOI: <https://doi.org/10.1016/j.ins.2019.11.045>

## Journal Pre-proof

Recommending investors for new startups by integrating network diffusion and investors' domain preference

Shuqi Xu, Qianming Zhang, Linyuan Lü, Manuel Sebastian Mariani

PII: S0020-0255(19)31090-4  
DOI: <https://doi.org/10.1016/j.ins.2019.11.045>  
Reference: INS 15037



To appear in: *Information Sciences*

Received date: 30 November 2018  
Revised date: 21 November 2019  
Accepted date: 25 November 2019

Please cite this article as: Shuqi Xu, Qianming Zhang, Linyuan Lü, Manuel Sebastian Mariani, Recommending investors for new startups by integrating network diffusion and investors' domain preference, *Information Sciences* (2019), doi: <https://doi.org/10.1016/j.ins.2019.11.045>

This is a PDF file of an article that has undergone enhancements after acceptance, such as the addition of a cover page and metadata, and formatting for readability, but it is not yet the definitive version of record. This version will undergo additional copyediting, typesetting and review before it is published in its final form, but we are providing this version to give early visibility of the article. Please note that, during the production process, errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

© 2019 Published by Elsevier Inc.

**Highlights**

- Network-based recommendation technique is a sharp tool to solve the problem of finding relevant investors for newcomer startups.
- Investors have their strong investment preferences in industry categories when selecting the invested startups.
- Recommending investor for new startups based on the tripartite network representations which include virtual connections between investors and tags is more efficient than the recommendation based on the straightforward network representation that only considers the direct relationship between investors, companies, and tags, as well as the existing collaborative filtering and matrix factorization algorithms.
- Among the three tripartite network representations, i.e., tag-company-investor, company-investor-tag, and company-tag-investor, the diffusion algorithm Probs on the weighted company-investor-tag network performs best.

# Recommending investors for new startups by integrating network diffusion and investors' domain preference

Shuqi Xu<sup>a</sup>, Qianming Zhang<sup>b</sup>, Linyuan Lü<sup>a,c,\*</sup>, Manuel Sebastian Mariani<sup>a,d</sup>

<sup>a</sup>*Institute of Fundamental and Frontier Science, University of Electronic Science and Technology of China, Chengdu, China*

<sup>b</sup>*Complex Lab, Big Data Research Center, University of Electronic Science and Technology of China, Chengdu, China*

<sup>c</sup>*Alibaba Research Center for Complexity Sciences, Hangzhou Normal University, Hangzhou, China*

<sup>d</sup>*URPP Social Networks, Universität Zürich, Zürich, Switzerland*

## Abstract

Over the past decade, many startups have sprung up, which create a huge demand for financial support from venture investors. However, due to the information asymmetry between investors and companies, the financing process is usually challenging and time-consuming, especially for the startups that have not yet obtained any investment. Because of this, effective data-driven techniques to automatically match startups with potentially relevant investors would be highly desirable. Here, we analyze 34,469 valid investment events collected from *www.itjuzi.com* and consider the cold-start problem of recommending investors for new startups. We address this problem by constructing different tripartite network representations of the data where nodes represent investors, companies, and companies' domains. First, we find that investors have strong domain preferences when investing, which motivates us to introduce virtual links between investors and investment domains in the tripartite network construction. Our analysis of the recommendation performance of diffusion-based algorithms applied to various network representations indicates that prospective investors for new startups are effectively revealed by integrating network diffusion processes with investors' domain preference.

**Keywords:** Diffusion model, Recommender systems, Venture investment, Cold start problem, Tripartite network

## 1. Introduction

The rapid development of the Internet brings the information overload problem: people usually receive too much information about a given issue, which greatly impairs the efficiency of their decision-making process [11]. Automated information-filtering tools such as ranking algorithms [23] and recommender systems [26] provide us with effective solutions to this problem. In particular, recommender systems can exploit data on users and their past preferences to predict their possible future interests, thus have been widely applied by various online platforms, including e-commerce websites [25] and online social networks [40]. Previous works [34] have shown that capable recommender systems can largely increase not only economic benefits but also customer loyalty.

In several financial investment scenarios, including stock investment [14], real estates investment [15], crowdfunding project [3], and portfolio management [30], recommender systems have received increasing attention [50]. Less attention has been devoted to the design and applications of recommendation systems in the domain of venture investment, an emerging investment type that aims at offering seed funding to startup companies. Extant studies on the topic [37, 48] have focused on helping VC firms to find investee companies, while effective methods to support new startups in their search for investors have not yet gained attention from scholars. For a startup, early-stage fundings constitute key support, yet obtaining them is challenging [32]. Finding a suitable investor who is interested in the startup's business scope usually requires long-term research, which is especially difficult for new startups due to their inexperience [8]. For this reason, an investor-filtering system aimed at new startups can be extremely beneficial.

\*Email address: linyuan.lv@gmail.com

Our main goal is to fill this gap by designing and validating recommendation system techniques to identify suitable investors for new startups. To this end, we analyze 34,469 investment events collected from [www.itjuzi.com](http://www.itjuzi.com) [1]. As we focus on startups that received no previous investments in the past, widely-studied recommendation techniques based on bipartite networks [49] are not applicable here since the new startups with no previous investors turn out to be isolated nodes in the investor-company bipartite network, making them not reachable by diffusion processes. Our problem can be classified as a *cold start problem* [35]: When a new actor enters the system, there is insufficient past information to provide him/her with a recommendation [26].

To overcome the cold-start problem, we resort to tagging techniques [16]. On the [www.itjuzi.com](http://www.itjuzi.com) platform, each startup is required to provide several tags that define its business scope. We found that the investors have a strong preference toward a small number of tags – in particular, in the majority of cases, an investor tends to invest in startups that feature her favorite tag. Hence, motivated by this finding and the key role played by the industry field in investment decision making [10, 28], we use tags as a key piece of information to generate recommendations.

By leveraging startups' tag information, we construct three different tripartite network representations of the investment system. The most natural tripartite representation is one where investors are connected to the startups they invested in, and startups are connected to their self-reported tags – we refer to this representation as the tag-company-investor (TCI) representation. This representation is natural because it is directly based on the collected data (companies' self-reported tags and investor-company investment events), and it is commonly used in the existing studies on online social networks [36, 46, 47]. This kind of tripartite networks only considers the connections that are already in the data, and the recommendation is made by applying physical diffusion processes such as probabilistic spreading [49] and heat transfer [45] (see details in Section 4.1) to the available network. However, we find that if our goal is to produce recommendation of investors for new startups, the natural TCI representation is suboptimal. The reason is that a new startup has not yet direct connections with investors, which implies that in the TCI representation, relatively long network paths are needed for a diffusion process to travel from a target startup to a prospective investor. Moreover, our empirical analysis indicates that investors tend to concentrate their investments toward startups with a small number of preferred tags. This property suggests that for our problem, investor-tag connections are potentially more informative than investor-company connections.

Motivated by these observations, we introduce virtual links between investors and tags, and construct two new tripartite network representations: the company-tag-investor (CTI) and the company-investor-tag (CIT) network representations. On each of the three tripartite network representations considered here (TCI, CTI, and CIT), we apply diffusion-based recommendation algorithms. We find that diffusion algorithms based on networks with virtual links (i.e., the company-tag-investor and company-investor-tag network) achieve substantially better recommendation accuracy compared to those based on the natural tag-company-investor representation.

Our main focus on diffusion-based techniques is motivated by existing works [26, 44] that indicate that this class of techniques is well-suited to scenarios where the input data are binary (without ratings) and only limited information about the target items or users is available. At the same time, we do not narrow our focus to this class of algorithms, alternative techniques from existing studies, including neighbor-based collaborative filtering [33] and matrix factorization based on Bayesian personalized ranking [13] (see Section 4.2) are also included in the analysis. The proposed diffusion-based algorithms turn out to outperform these existing algorithms that are not based on a tripartite network representation.

The main contribution of this paper is twofold. First, it brings network-based recommendation techniques to the problem of finding relevant investors for newcomer startups. Second, to address the problem, it introduces virtual links between investors and tags to build the network representation. The possibility to introduce these links has not been exploited by previous studies on recommendation algorithms based on tripartite networks [36, 46, 47], yet the virtual links turn out to be vital to achieving a good recommendation performance. The new representation achieves indeed improved performance over recommendation techniques based on the straightforward TCI representation and existing collaborative filtering [33] and matrix factorization algorithms [13]. Our findings reveal a powerful method for startups to find their first investor.

The paper is structured as follows: In Section 2, we review related works. In Section 3, we describe the dataset, analyze the tag preference of investors, and propose the construction of three tripartite network representations. In Section 4, we introduce the recommendation algorithms employed in the tripartite networks, the baseline algorithms, and the evaluation metrics. In Section 5, we present our results and discuss them. Finally, Section 6 concludes the paper and point out several open research directions.

## 2. Related work

### 2.1. Recommendation methods

Collaborative filtering (CF) is one of the most widely-used family of recommendation techniques. It can be further grouped into memory-based CF and model-based CF [39]. In online systems where users rate items, memory-based CF approaches use users' past rating data to compute the similarity between users or items and produce a prediction for the target user based on the similar users to her or similar items to those she already collected. Model-based CF methods leverage historical data to learn a predictive model. Well-known model-based CF techniques include Bayesian networks [29], clustering models [6], latent semantic models [19], among others. The most critical component of CF methods is the measurement of the similarity between pairs of users or items, which are known to be vulnerable against data sparsity and cold-start problem [2].

Another common class of techniques is content-based filtering, which utilizes the description of items and the profile of users' preference to find the items that best match the items previously selected by the user. The recommendation process includes representing the items' features, learning users' profile, and filtering. There are also hybrid approaches which combine collaborative with content-based methods or with different variants [4].

While most recommender systems act on data with ratings, positive-only data (e.g., click-through data, browsing history, investment history) without ratings are of interest to several physics-rooted recommendation methods. In this context, the input data are represented as a bipartite network (e.g., a user-item network [49]) or a tripartite network (e.g., a user-item-tag network [47]) based on the actual links in the dataset. The recommendation can be obtained by employing classical physics processes such as diffusion (like in the Probs algorithm [49]) and heat transfer (like in the Heats algorithm [45]) on the network. We refer to [44] for a review of network-based recommendation algorithms.

### 2.2. Recommender systems in investment domains

Applying recommender systems to financial investment problems has received growing attention [50]. It is generally considered as a challenging task because of the strict expectations from the information-seeker [50]. Fields of application for recommender systems in finance-related domains include stock investment [14], real estate investment [15], crowdfunding projects [3], and portfolio management [30], among others. As for venture investment, there have been relatively few scientific publications on the topic of recommendation. Stone et al. [37] used collaborative filtering to recommend relevant investment opportunities to venture capital (VC) firms. They reported that this class of activities is characterized by extremely sparse data, and the number of invested companies for VC firms follows a power-law distribution, which points out the existence of very active VC firms. Zhao et al. [48] proposed 5 risk-aware startup selection methods and ranking algorithms to predict VC firms new investments. Existing studies largely concentrated on helping VC firms to find investee companies, while there is a lack of research and effective methods to support new startups searching for investors. Bringing network-based recommendation techniques to the problem of recommending investors to startups is one of the main contributions of our paper.

### 2.3. Cold-start problems

In cold-start problems, because of insufficient past information, it is hard to infer users' preferences or items' potentially relevant users. A simple approach to mitigate the cold-start problem is to recommend the most popular objects based on historical data. But this strategy can only provide a uniform recommendation to all users. More diverse outcomes can be obtained by gathering rapidly additional information on users' preferences. This can be achieved by actively eliciting the user to make more informative choices [31], by integrating information from other user activities [5], or by using hybrid techniques to combine recommendations obtained by different methods [24]. As an alternative to these techniques, one can leverage the additional information to construct a network and run a standard diffusion-based algorithm on it. Following this idea, Zhang et al. [46] proposed a diffusion-based recommendation algorithm which considers social tags as a bridge connecting users and objects. They indicated that tags can effectively build up relations between existing objects and new ones, thereby providing solid recommendations for new objects. Deng et al. [9] introduced the Social Mass Diffusion (SMD) method based on a mass diffusion process in the combined network of users social network and user-item bipartite network. They showed that the SMD can generate more personalized recommendations for new users than the global ranking based on popularity.

investor	company	company's tags	time
IDG Capital	Tencent	social network	2000.04.01
		comprehensive social communication	
		comprehensive financial service	
		comprehensive game service	
		SEO/SEM	
Google	Baidu	platform	2004.06.01
		enterprise service	
		comprehensive enterprise service	
		search engines	
		local comprehensive life	

Table 1: Two examples of investment events that involve well-known investors and companies. Notice that all the reported entries are recorded in Chinese in the original dataset.

Besides, scholars introduced a number of strategies based on matrix factorization algorithm. Gantner et al. [13] proposed a method by an extension of matrix factorization optimized for Bayesian Personalized Ranking. They leveraged the mapping functions to compute adequate latent feature representations for new entities from their attributes. Kula [22] estimated feature embeddings by factorizing the collaborative interaction matrix. In his approach, new users or items can be represented in terms of combinations of metadata features that have been estimated from the training set. Fernández-Tobías et al. [12] analyzed several solutions to the new user problem in collaborative filtering based on users' personality information (using 5 factors to describe an individual's personality: openness, conscientiousness, extraversion, agreeableness and neuroticism, as assessed by self-descriptive sentences or adjectives), including personality-based matrix factorization, personality-based active learning, and personality-based cross-domain recommendation.

To summarize, the key element to address the cold-start problem is the acquired information about the users preference obtained either by directly asking for auxiliary information or by actively collecting it when available in the studied platforms. But requiring the targets to provide detailed individual information is not practicable in most cases. Therefore, the majority of approaches can only mitigate the cold-start problem when a target user has at least a limited historical activity, whereas recommendations for new users without prior activity must draw support from implicit preferences or features. In this paper, we focus on the cold-start problem for startups without prior activity, and we leverage the startups' self-reported tags to infer investors' preferences.

### 3. Data and network construction

#### 3.1. Dataset description

We collected 45,943 investment events from *www.itjuzi.com* [1]. Each event in the data includes the information of one investor, one investee company, the company's tags, and the event time (with the temporal resolution of one day). We note that joint investments that involve several investors appear as multiple events. We removed from the data all the events where the investor information is missing, which leaves us with 34,469 investment events ranging from Dec. 1<sup>st</sup>, 1999 to Apr. 28<sup>th</sup>, 2017. There are 5,588 investors, 14,887 companies, and 1,080 tags involved in these events. Table 1 shows two examples.

On average, an investor invests in 6.2 companies – the blue squares in Figure 1 (a) represent the distribution of the number of companies per investor. On the other hand, on average, a company has 2.3 investors – the red dots in Figure 1 (a) represent the distribution of the number of investors per company. 52% of companies have only *one investor*, which is a manifestation of the sparse nature of investment data. Such a sparsity is reasonable since investment events are not as frequent as online activity events (such as watching movies or making new friends in online platforms) [48], and the final investment decision is usually time-consuming [8]. As for tags, each company reported 4.9 tags, on average – Figure 1 (b) illustrates the distribution of the number of tags per company.

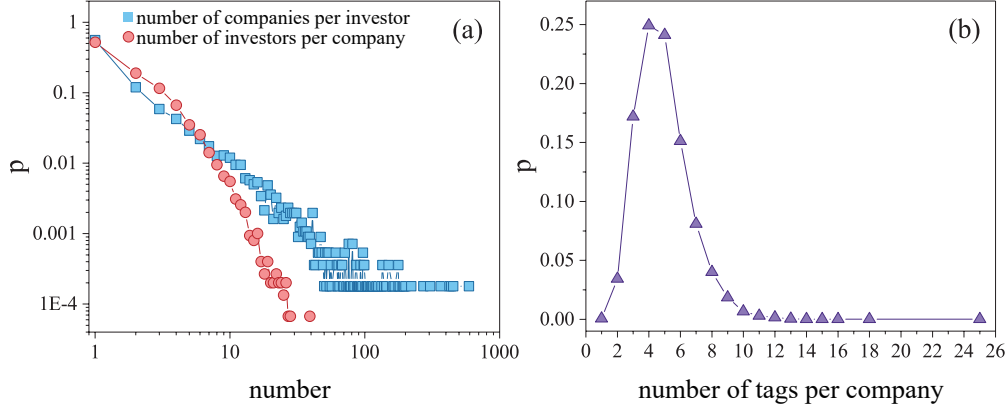


Figure 1: Statistical properties of the investment dataset. In panel (a), the blue squares represent the distribution of the number of companies per investor, whereas the red dots represent the distribution of the number of investors per company. Panel (b) shows the distribution of the number of tags per company.

### 3.2. The strong domain preference of investors

Our first main result is that investors tend to invest in specific tags (which correspond to specific domains). Evidence for the strong preference of investors toward specific tags is reported in Figure 2. First, for each investor  $I_x$ , we determine her most frequent investing tag (denoted as  $T_i^*(I_x)$ ) among those that belong to the companies that  $I_x$  has invested in. Subsequently, we determine the fraction of  $I_x$ 's investments toward startups that feature tag  $T_i^*(I_x)$  – this fraction represents the strength of  $I_x$ 's preference toward its preferred tag, and we denote it as  $P(I_x)$ . The distribution of  $P(I_x)$  for all investors is shown in Figure 2 (a). For each investor, on average, 45% of the companies it invested in feature its preferred tag. This indicates that many investors narrow down their focus to a limited number of fields, and their preferred field dominates their investment activity. To further illustrate this property, in Figure 2 (b), we assess the similarity between companies in which an investor has invested. More specifically, for each investor, we build a company similarity network where two companies are connected if they share at least one common tag. We consider the fraction of companies that belong to the giant component of the network as the measure of company similarity for investor  $I_x$ , and we denote it as  $P'(I_x)$ . The distribution of  $P'(I_x)$  for all investors is shown in Figure 2 (b), and the average value is as high as 52%.

Overall, Figure 2 reveals investors prefer to focus on a few industry categories. The reason might be that investors are cautious in investing in unfamiliar fields to avoid risk [28], which results in a low diversity in their investment activity.

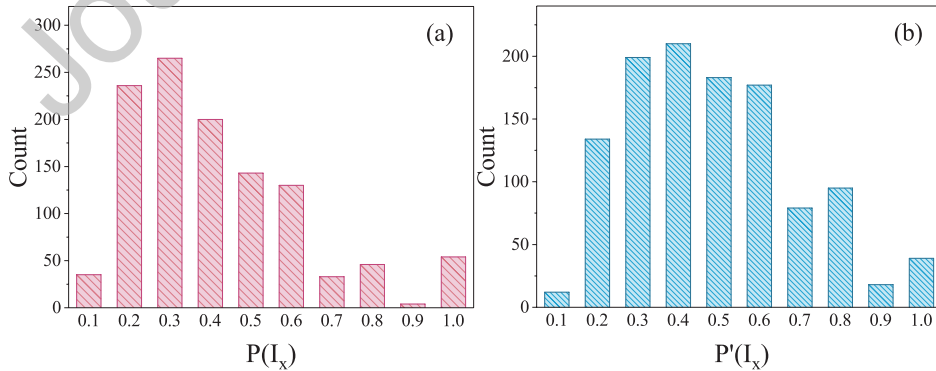


Figure 2: The strong domain preference of investors. Investors invested in less than five companies are not taken into account since their preferences are relatively difficult to capture. Panel (a) shows the distribution of the probability that one investor invests in a company with the most preferred tag of this investor. Panel (b) shows the distribution of the fraction of companies that belong to the giant component of the similarity network of the companies in which the investor has invested.



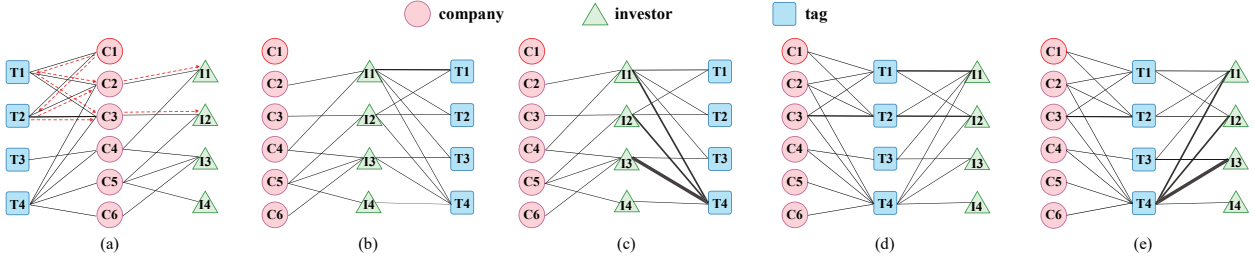


Figure 3: An illustration of different possible tripartite network representations. The thickness of each line is proportional to its weight. The weight of the link between a given investor and a given tag indicates the number of companies with that tag the investor has invested in. (a) Traditional tag-company-investor (TCI) network. (b) Unweighted company-investor-tag (CIT) network. (c) Weighted company-investor-tag (CIT) network. (d) Unweighted company-tag-investor (CTI) network. (e) Weighted company-tag-investor (CTI) network.

### 3.3. Tripartite network construction

Investment relationships can be modeled by means of an investor-company bipartite network, which represents interactions between investors and startup companies. However, in our problem of interest, the target startups have no past connections with any investors. Therefore, none of the existing recommendation techniques based on bipartite network representations are relevant to the problem. A natural network representation for our problem is the tag-company-investor (TCI) tripartite network, which simply integrates the investor-company bipartite network and the company-tag bipartite network, as shown in Figure 3 (a). An analogous representation is commonly used in existing studies on online social networks [36, 46, 47]. Differently from online systems where users can freely assign tags to their collected items, in investment systems, companies' tags are self-reported and chosen from a tag library.

Formally, we denote the TCI network representation as  $\mathcal{G}(T, C, I)$ , where  $T, C, I$  represent the set of tags, companies and investors, respectively. In this network, link  $e(I_x, C_i)$  means that investor  $I_x$  invested in company  $C_i$ , whereas link  $e(C_i, T_a)$  means that company  $C_i$  has the tag  $T_a$ . Considering the problem where a given target startup –  $C_1$  in Fig. 3 – needs a recommendation, yet it has no direct connections to investors, we are not able to find any related investors unless following the red dashed arrows to identify the investors who have invested in companies with common tags to the target company. This is also the way how tags work in recommending investors to startups through the TCI network representation. So tags act as the bridge between startups without investors and startups that already received investments.

However, this type of tripartite network has one marked disadvantage. A popular tag which is connected to many companies, like  $T_4$  in Figure 3 (a), will connect the target startup with many companies, and among them, unrelated companies may bring uninterested investors into consideration. Therefore, a startup that owns some popular tags, like  $C_2$ , will find many noisy investors following network paths. Besides, we found that many investors prefer to invest in some specific field with certain tags (Section 3.2), which cannot be leverage by the TCI representation, as tags and investors are on the opposite sides of the tripartite network.

To overcome the limitations of the TCI network representation, we construct two new tripartite networks, i.e., the company-tag-investor (CTI,  $\mathcal{G}(C, T, I)$ ) and the company-investor-tag (CIT,  $\mathcal{G}(C, I, T)$ ) network. These two representations include virtual connections between investors and tags. We create these virtual links between investors and tags as follows: if investor  $I_x$  invested in company  $C_a$ , and  $C_a$  has tags  $T_i$  and  $T_j$ , then links  $e(I_x, T_i)$  and  $e(I_x, T_j)$  are created. Moreover, we can also construct the corresponding weighted representations by introducing weighted links: for example, weight  $w_{xi} = 2$  means that investor  $I_x$  invested in two companies with tag  $T_i$ . As shown in Figure 3 (b)-(e), in the constructed networks, there are direct links between investors and tags which can be either weighted or unweighted.

## 4. Methods

### 4.1. Recommendation algorithms on tripartite networks

Based on tripartite networks, the probabilistic spreading (Probs) recommendation algorithm was first proposed for bipartite networks [49] and later proved to be effective for tripartite networks [46, 47]. The Probs algorithm assigns

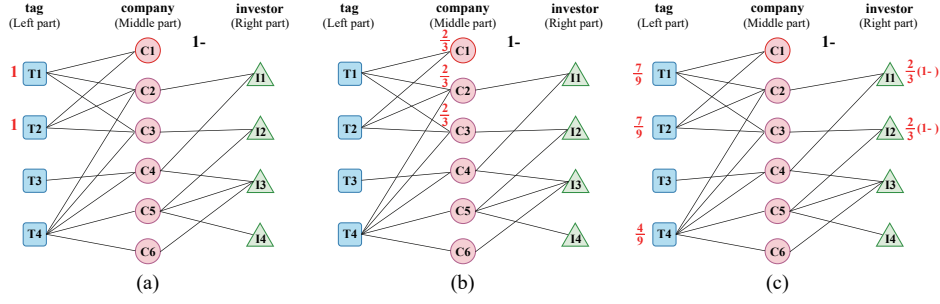


Figure 4: Implementing Probs algorithm on tag-company-investor tripartite network. In (a), given the target company  $C_1$ , we assign  $C_1$ 's tags, i.e.,  $T_1$  and  $T_2$ , one unit resource. In (b), all nodes redistribute their resources to their neighbors, and at the same time, all nodes update their resource values following Eq. (2). In (c), all nodes redistribute their resources to their neighbors once again.

an initial resource to specific nodes, then redistributes this resource in a similar way to a random-walk process. In the following, to generalize the representations in Fig 3, we consider a generic L-M-R tripartite network representation and denote the three classes of nodes as Left nodes (L), Middle nodes (M) and Right nodes (R), respectively. Suppose that initially, each node possesses a kind of resource, which are denoted as  $f(L_x)$ ,  $f(M_y)$ ,  $f(R_z)$ , namely the resource possessed by left node  $L_x$ , middle node  $M_y$ , and right node  $R_z$ , respectively. Subsequently, each node distributes its resource to all its neighbors. The new resource values for Left, Middle and Right nodes are determined by the equations,

$$\begin{cases} f'(L_x) = \sum_{M_y \in M} \frac{f(M_y)A(L_x, M_y)}{k_{M_y \rightarrow L}}, \\ f'(R_z) = \sum_{M_y \in M} \frac{f(M_y)A(R_z, M_y)}{k_{M_y \rightarrow R}}, \\ f'(M_y) = \sum_{L_x \in L} \frac{f(L_x)A(L_x, M_y)}{k_{L_x}} + \sum_{R_z \in R} \frac{f(R_z)A(R_z, M_y)}{k_{R_z}}, \end{cases} \quad (1)$$

where  $A$  is the adjacency matrix,  $A(L_x, M_y) = 1$  if  $L_x$  and  $M_y$  are connected.  $k_{M_y \rightarrow L}$  is the number of neighboring left nodes for  $M_y$ , while  $k_{M_y \rightarrow R}$  is the number of connected right nodes for  $M_y$ .  $k_{L_x}$  and  $k_{R_z}$  mean the degree of  $L_x$  and  $R_z$ , respectively. Note that if the links are weighted,  $k$  represents the sum of weights on the related links.

Considering that the middle nodes have two kinds of neighbors whose importance may be different, we introduce a tunable parameter  $\lambda \in [0, 1]$  as in [47]. During the process of distribution, each middle node distributes a ratio  $\lambda$  of its resources to the neighboring left nodes, and the rest  $1 - \lambda$  resources are distributed to the neighboring right nodes. We calculate the new resource through

$$\begin{cases} f'(L_x) = \lambda \sum_{M_y \in M} \frac{f(M_y)A(L_x, M_y)}{k_{M_y \rightarrow L}}, \\ f'(R_z) = (1 - \lambda) \sum_{M_y \in M} \frac{f(M_y)A(R_z, M_y)}{k_{M_y \rightarrow R}}, \\ f'(M_y) = \sum_{L_x \in L} \frac{f(L_x)A(L_x, M_y)}{k_{L_x}} + \sum_{R_z \in R} \frac{f(R_z)A(R_z, M_y)}{k_{R_z}}. \end{cases} \quad (2)$$

Given a target startup which needs recommendation and one of the tripartite network representations in Figure 3, we start the iterative process described above by assigning one unit resource to all the tags of the target startup. Subsequently, the resource spreads along the network edges. A visual representation of the 2-step Probs process is given in Figure 4. In this example, investors can only receive resource and update their score (new resource value) at even steps (step two, four, six, ...). By contrast, in the company-investor-tag and company-tag-investor representations, investors can obtain new resources and update their scores at odd steps (step one, three, five, ...). To perform a clear comparison between different tripartite networks with different diffusion steps of the algorithms, we replace *step* with *reach* by referring to the state of the resource after the resource arrives at the investors for the  $i$ -th time as the  $i$ -reach of the process. In the TCI representation as an example, 1-reach represents the 2<sup>nd</sup> step, 2-reach represents the 4<sup>th</sup> step, and so forth. In the CIT representation, 1-reach represents the 1<sup>st</sup> step, 2-reach means 3<sup>rd</sup> step, and so forth.

For the sake of completeness, we also consider an alternative diffusion process on the tripartite network representations introduced above: the heat spreading (Heats) algorithm [45], which employs a process analogous to heat diffusion. We denote by  $h(L_i)$  the level of “heat” of node  $L_i$ . The new level of “heat” of each node after diffusion process are calculated following these equations,

$$\begin{cases} h'(L_x) = \lambda \sum_{M_y \in M} \frac{h(M_y)A(L_x, M_y)}{k_{L_x}}, \\ h'(R_z) = (1 - \lambda) \sum_{M_y \in M} \frac{h(M_y)A(R_z, M_y)}{k_{R_z}}, \\ h'(M_y) = \sum_{L_x \in L} \frac{h(L_x)A(L_x, M_y)}{k_{M_y \rightarrow L}} + \sum_{R_z \in R} \frac{h(R_z)A(R_z, M_y)}{k_{M_y \rightarrow R}}. \end{cases} \quad (3)$$

Similar to Probs, we can get  $i$ -reach Heats.

For Probs or Heats, each investor can obtain the value of “resource” (or “heat”) after  $i$ -reach diffusion process ( $i = 1, 2, 3, \dots$ ). The investors with higher values will be the top choice to build the recommendation list. The working process of our approach is shown in Figure 5.

#### 4.2. Benchmark algorithms

The techniques introduced above build on diffusion processes on tripartite network representations. Besides assessing their relative performance, it is essential to compare their performance against simpler algorithms that do not require an iterative process to compute the recommendation scores, and against state-of-the-art recommendation techniques that are potentially relevant to our recommendation problem. In total, we consider five algorithms as benchmark algorithms: three of them are simple metrics, whereas two of them are more sophisticated, state-of-the-art methods. We provide below the details of the benchmark algorithms.

*Popularity-based metric.* We measure the number of investments that each investor has previously made – i.e., the investors’ degree in the investor-company bipartite network as the recommendation score. The more investments an investor has previously made, the higher its ranking position in the recommendation list for new startups that did not yet receive any investment.

*Tag-voting metric.* We consider simple metrics that leverage tag information. For a target startup  $C_a$  with tags  $\{T_1, \dots, T_k\}$ , if investor  $I_x$  invested in  $n_i$  companies with tag  $T_i$ , the score of  $I_x$  for  $C_a$  is defined as

$$s(I_x, C_a) = \sum_{i=1}^k n_i \quad (4)$$

We refer to the recommendation of investors based on this score (the investors with the largest score are recommended to the target startup) as the tag-voting method.

*Normalized tag-voting metric.* To prevent the potential bias toward popular tags of the tag-voting metric, we further consider a normalized metric defined as

$$s(I_x, C_a) = \sum_{i=1}^k \frac{n_i}{N_i} \quad (5)$$

where  $N_i$  is the total number of investments received by companies with tag  $T_i$ . We refer to the recommendation based on this metric (again, the investors with the largest score are recommended to the target startup) as the normalized tag-voting method.

*Neighbor-based Collaborative filtering.* Collaborative filtering technology is by far one of the most widely used and successful method [21]. A large number of studies has proved its high accuracy in various systems, such as E-commerce platforms [25], digital research libraries [41], online social systems [43], and so on. Here we also tested an extend collaborative filtering method as a benchmark. Using tag information, one can measure the similarity between

two companies, and leverage the investment history of similar companies to generate a recommendation. According to [33], the score of investor  $I_x$  for startup  $C_a$  is defined as

$$s(I_x, C_a) = \frac{\sum_b^{N_c} \text{sim}(C_a, C_b) A(I_x, C_b)}{\sum_b^{N_c} \text{sim}(C_a, C_b)} \quad (6)$$

where  $N_c$  is the number of companies in the dataset,  $\text{sim}(C_a, C_b)$  is the similarity between company  $C_a$  and  $C_b$ ,  $A(I_x, C_b) = 1$  if investor  $I_x$  has invested in company  $C_b$ . Following [33], in this paper, we leverage the cosine similarity [33] which is given by

$$\text{sim}(C_a, C_b) = \frac{|T(a) \cap T(b)|}{\sqrt{|T(a)| \times |T(b)|}} \quad (7)$$

where  $T(a)$  refer to the tags belong to company  $C_a$ ,  $|T(a)|$  is the size of  $T(a)$ ,  $T(a) \cap T(b)$  is the common tags for  $C_a$  and  $C_b$ . We name this metric as neighbor-based CF.

*Matrix factorization based on Bayesian Personalized Ranking.* Matrix factorization with Bayesian Personalized Ranking learning (BPR) framework has been applied to various recommender systems and generally considered as a powerful recommendation method in implicit or positive-only feedback dataset [7, 20], which fit the characteristic of our collected data – investor behavior is in one-class (investing) form. We consider the combination of the matrix factorization based on BPR and feature mapping. Following Gantner et al.’s work [13], we used a matrix factorization model based on the BPR framework to train the latent-feature factor of investors and companies. Then the latent-feature vector of a given new startup can be approximated based on those of similar companies. With this estimation, we can compute a score for the new startup  $C_a$  with an investor  $I_x$  through

$$s(I_x, C_a) = \langle W(I_x), \phi(C_a) \rangle \quad (8)$$

where  $W(I_x)$  is the trained factor for  $I_x$ ,  $\phi(C_a)$  is  $C_a$ ’s estimated latent factor. The similarity between companies is determined through Equation (7). And we employ the weighted  $k$ -nearest-neighbor (kNN) regression [42] to estimate the factor of the new startup. For example,  $C_a$ ’s latent factor is estimated through the equation as in [13]

$$\phi(C_a) = \frac{\sum_{b \in N_k(C_a)} \text{sim}(C_a, C_b) \phi(C_b)}{\sum_{b \in N_k(C_a)} \text{sim}(C_a, C_b)} \quad (9)$$

where  $N_k(C_a)$  represents the set that comprises the  $k$  most similar companies to  $C_a$  and  $k$  is set to 6 as it is a suitable value on the training set. The factor number is equal to 30 using hyperparameters that typically yields satisfactory results in non-cold start evaluations in our dataset. We also tested different numbers of factors (10, 20, 30, 40, 50, 60), and results are qualitatively similar. We refer to this method as BPR-MF method.

#### 4.3. Evaluation

To evaluate the proposed recommendation methods, we divide the dataset into two parts, the training set and the testing set. The training set includes 31,039 investment events that happened before Sept. 6, 2016 (i.e., the 90% earliest events in the data). Based on it, we construct the different tripartite network representations described above. The remaining 3,430 investment events (i.e., the 10% latest events in the dataset) compose the testing set. For the validation of information-filtering predictive techniques, there is no unique or universal criterion to choose the relative size or temporal duration of training and testing set. The 90%/10% proportion adopted here a standard choice in the evaluation of recommender systems. Yet, we also tested different sizes (85%/15% and 95%/5%) – the obtained results are in qualitative agreement with those obtained with the 90%/10% splitting, and they are discussed in Appendix A.8 and A.9. We stress that our aim is to recommend investors to new startups; as a consequence, only the 1,096 companies that obtained their first investment within the testing set and their corresponding investors are considered in the evaluation – we refer to them as *target startups*.

2019/9/28

(1).svg

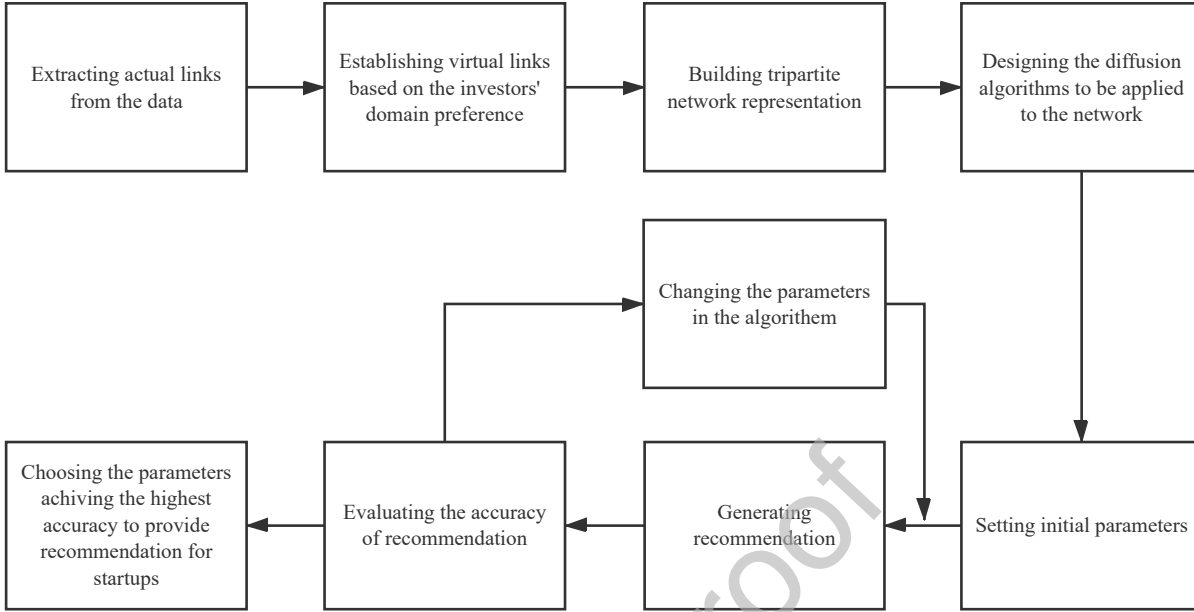


Figure 5: The working flow of the proposed recommendation process.

*Ranking Score (RS)* [49] is used to evaluate the various recommendation methods. For a target startup, we refer to the investors who actually invested in the target startup within the testing set as the *relevant investors* to that startup and refer to all the other investors as *irrelevant investors*. The *RS* measures the relative ranking of the relevant investors in the target startups' recommendation list: when there are  $o$  investors that can potentially be recommended, a relevant investor with ranking  $r$  achieves the relative ranking  $r/o$ , which is equal to her ranking score. By averaging over all target startups in the testing set and their relevant investors, we obtain the mean *RS*: the smaller the mean ranking score, the higher the algorithm's accuracy. We stress that to evaluate the recommendation methods' performance in the cold-start problem studied here, the ranking score is a better evaluation metric than *Precision* and *Recall* because the target startups typically have only few investors.

We also measure the *AUC* (Area under receiver operator curve) [18] as an alternative performance metric. The *AUC* aims at assessing the method's ability to distinguish relevant from irrelevant investors. For a target company, its *AUC* can be approximated as the probability that, when choosing at random one investor from the relevant investors and another investor from the irrelevant investors, and the relevant investor's score is higher than the irrelevant investor's score. We get the mean *AUC* by averaging over all target companies in the test set; the higher the *AUC*, the better the method's performance.

## 5. Results

Figure 6 compares the recommendation results for the Probs algorithm for different network representations. Figure 6 (a) and (b) show the *RS* values for 2-reach Probs diffusion and 3-reach Probs diffusion, respectively, as a function of  $\lambda$ , which represents the proportion of resources distributed from central to left and right nodes. Generally, the proposed tripartite network representation that leverages virtual links between investors and tags (i.e., company-investor-tag and company-tag-investor tripartite networks) outperforms the natural tag-company-investor representation. Besides, algorithms based on weighted network representations outperform those based on unweighted representations (see Figure 6).

Comparing Figures 6 (a) and (b), we find that the average *RS* depends on the reach of the considered diffusion process. To identify the best overall method, we obtain the best *RS* values for all the tripartite networks at 2-reach diffusion, 3-reach diffusion, ..., 7-reach diffusion respectively, and show them in Figure 6 (c). We conclude that

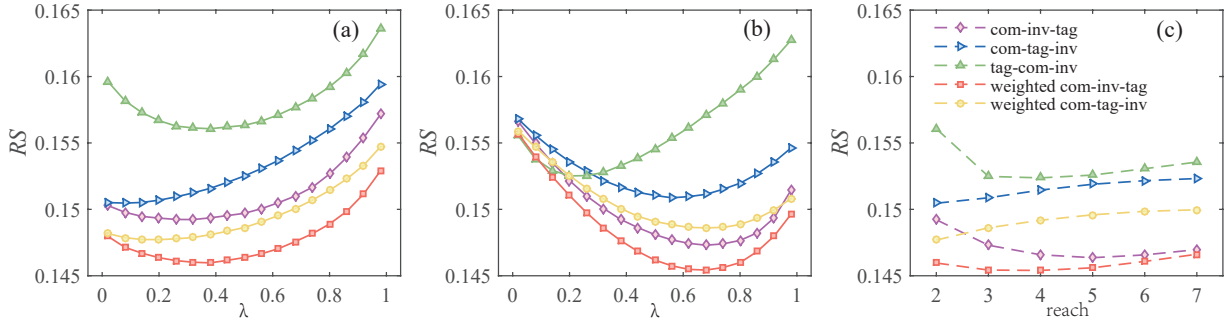


Figure 6: The comparison of Probs recommendation algorithms applied to different tripartite networks. Each line corresponds to a specific network representation. (a) and (b) show the average  $RS$  values for 2-reach Probs and 3-reach Probs, respectively, as a function of the algorithm parameter  $\lambda$ . (c) shows the average  $RS$  value attained by the Probs algorithm with optimal  $\lambda$ , as a function of the reach of the diffusion process. Note that in the legends, com and inv stand for investor and company, respectively.

both the optimal  $RS$  curve (see Figure 6 (c)) and the optimal  $RS$  value (see Table 2) are obtained from the weighted company-investor-tag network.

Network	Com-inv-tag		Com-tag-inv		Tag-com-inv
	weighted inv-tag	unweighted inv-tag	weighted inv-tag	unweighted inv-tag	unweighted
Probs 1-reach	0.21094	0.22334	0.21094	0.22334	0.21478
Probs 2-reach	0.14596	0.14922	0.14769	0.15046	0.15607
	$\lambda^*=0.36$	$\lambda^*=0.28$	$\lambda^*=0.18$	$\lambda^*=0.10$	$\lambda^*=0.38$
Probs 3-reach	0.14543	0.14730	0.14860	0.15087	0.15248
	$\lambda^*=0.68$	$\lambda^*=0.68$	$\lambda^*=0.66$	$\lambda^*=0.56$	$\lambda^*=0.22$
Probs reach*	<b>0.14540</b>	0.14636	0.14769	0.15046	0.15237
	$\lambda^*=0.82$	$\lambda^*=0.92$	$\lambda^*=0.18$	$\lambda^*=0.10$	$\lambda^*=0.14$
	<b>reach*=4</b>	reach*=5	reach*=2	reach*=2	reach*=4
Heats reach*	0.25612	0.27259	0.27073	0.27880	0.25631
	$\lambda^*=0.56$	$\lambda^*=0.44$	$\lambda^*=0.34$	$\lambda^*=0.32$	$\lambda^*=0.60$
	reach*=2	reach*=2	reach*=2	reach*=2	reach*=2
benchmarks	popularity-based method=0.16507 tag-voting = 0.20781 normalized tag-voting = 0.21094 neighbor-based CF = 0.20791 BPR-MF = 0.35971				

Table 2: Average  $RS$  of the considered recommendation algorithms for three tripartite network representations. For each value of reach, we provide the optimal  $\lambda$  value together with the corresponding average  $RS$ . We also provide the performance by the Heats algorithm and three local benchmark metrics. For each network representation, the Probs algorithm substantially outperforms the Heats algorithm. It also outperforms the benchmark metrics.

It is essential to compare the performance of the Probs method with both the Heats algorithm and the benchmark algorithms introduced in Section 4.2. We find that the optimal performances achieved by the Heats algorithm (i.e., the performance achieved by selecting the optimal reach of the algorithm, for each network representation) are substantially worse than those achieved by non-optimized Probs algorithm (see Table 2). For instance, the average  $RS$  of Heats algorithm achieves 0.25612 with 2-reach on the weighted company-investor-tag tripartite network, as compared to 0.14596 achieved by Probs algorithm with 2-reach on the same tripartite network. Similarly, the average  $RS$  of the three simple benchmark metrics introduced above (popularity-based, tag-voting, and normalized tag-voting metric) is larger than that obtained by the Probs algorithm with more than one reach (see Table 2 and Figure 7 (a)). For example,

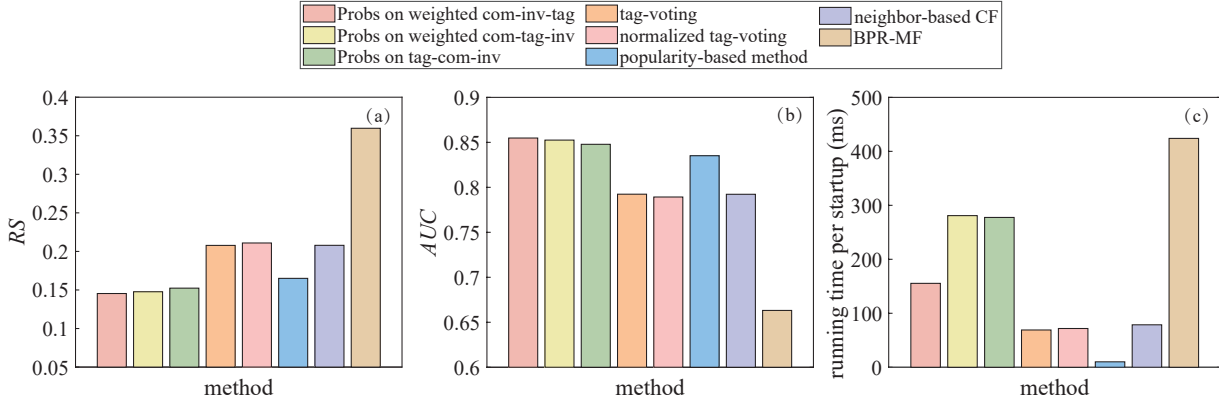


Figure 7: (a) Average  $RS$  of the Probs algorithm for three tripartite network representations and five benchmark methods. The applied  $\lambda$  and reach for the Probs algorithm are those achieve the optimal performance as shown in Table 2. The smaller the  $RS$ , the better the methods performance. (b) Average  $AUC$  of the Probs algorithm for three tripartite network representations and five benchmark methods. Noted that  $\lambda$  and reach are chosen based on the optimal performance in the  $RS$  evaluation. The higher the  $AUC$ , the better the methods performance. (c) The test times per startup for different methods. We test 2-reach Probs on the three network representations.

among these local metrics, the best-performing one (popularity-based method) achieves average  $RS$  equal to 0.16507, as opposed to 0.14596 achieved by Probs with two reach based on the weighted company-investor-tag tripartite network. As for the common-used collaborative filtering and advanced matrix factorization model, their performance is also inferior to that of the Probs algorithm on tripartite networks. This could be due to the sparsity of the investment data that may limit their effectiveness [38]. And the results also indicate that the tag information is more powerful when it is connected to the investors' preference than when it is used to measure the similarity between companies.

To further verify the stability of the results, we repeat the predictive analysis for different proportions of the training and testing set sizes (85%/15% and 95%/5%), and for another evaluation metric ( $AUC$ ). It turns out that the relative performance of the methods is in agreement with the findings for the 90%/10% splitting (see Figure A.8 and A.9). The  $AUC$  results as shown in Figure 7 (b) also show the advantage of the accuracy of Probs on weighted com-inv-tag network representation, which are in agreement with the  $RS$  results.

We conclude that: (i) the Probs algorithm performs significantly better than the Heats algorithm in the recommendation of investors to new startups; (ii) the benchmark metrics cannot compete with the Probs algorithm on the tripartite network representations, which suggests that constructing the tripartite networks considered here are indispensable for an accurate recommendation. (iii) the Probs algorithm on the weighted com-inv-tag achieves the best recommendation performance.

As for run-time overhead (Figure 7 (c)), one can clearly see that 2-reach Probs in the weighted com-inv-tag network has a relatively fast testing time, while the other two network representations take longer times to generate the predictions. Overall, the computational time of all the considered metric is of the same order of magnitude.

## 6. Conclusion and discussion

The main contribution of this paper is to solve the investor recommendation problem by building tripartite network representations including virtual links in combination with the Probs diffusion algorithm. The new application scenario and the introduction of virtual links into diffusion-based recommendation are the main innovative features of our study. Obtaining financial support is a rigid demand for most startups; thus, to precisely recommend investors for new startups is valuable, yet difficult due to the cold-start problem. We started by constructing a natural tag-company-investor tripartite network representation. In consideration of its disadvantages and the significant role played by tag in investment decisions, we reconstructed two tripartite network representations by establishing virtual connections between investors and tags, which take full advantage of the investors' preference for tags. The obtained results demonstrate that our modification of the traditional tripartite network is effective to provide better recommendations, and the weighted company-investor-tag network achieves the best performance in accuracy. Besides, the Probs algorithm significantly outperforms the Heats algorithm and baseline methods.

Recommendation in financial investment domains – especially in the venture investing domain – is a new topic in the recommender systems literature. A few academic studies have worked on applying recommendation techniques to provide personalized recommendations for investors – whether individual investors or investment institutions – but neglected the startups’ demand. Our research paves the way to the extensive investigation of the investor recommendation problem for startups.

We identify three main directions for future improvements. First, we have not analyzed the temporal effects of related data. A fundamental assumption of our method is that investors tend to invest in companies with tags which are similar to those of their past investments. But in reality, the preference of investors can change over time. Incorporating the temporal information may improve the recommendation performance. Second, it is essential to test our new tripartite networks along with other recommendation methods on more investment datasets. These attempts will provide us a promising way to better inform investment decision making and, as a result, assist startups effectively. Third, we leveraged domain (tag) information to generate effective recommendations. Incorporating into the recommendation algorithm additional factors such as financial consideration, product and market characteristics, social relation information [17, 27, 28], might further improve the predictive performance.

### Acknowledgements

This work is supported by the National Natural Science Foundation of China (Grants nos. 11622538, 61673150, 61703074), the Zhejiang Provincial Natural Science Foundation of China (Grant no. LR16A050001), and the Science Strength Promotion Programme of UESTC. MSM also acknowledges support from the University of Zurich through the URPP Social Networks, the Swiss National Science Foundation Grant No. 200021-182659, the UESTC professor research start-up grant No. ZYGX2018KYQD215.

### Author contributions statement

Q.Z. and L.L. conceived the idea and designed the research. S.X. and Q.Z. collected and processed the data. S.X. performed the experiments. All authors analyzed and discussed the results. S.X., Q.Z. and M.S.M. wrote the manuscript. All authors reviewed the manuscript.

### References

#### References

- [1] <http://www.itjuzi.com/>; [accessed 1 September 2018].
- [2] Hyung Jun Ahn. A new similarity measure for collaborative filtering to alleviate the new user cold-starting problem. *Information Sciences*, 178(1):37–51, 2008.
- [3] Jisun An, Daniele Quercia, and Jon Crowcroft. Recommending investors for crowdfunding projects. In *Proceedings of the 23rd international conference on World wide web*, pages 261–270. ACM, 2014.
- [4] Robin Burke. Hybrid recommender systems: Survey and experiments. *User modeling and user-adapted interaction*, 12(4):331–370, 2002.
- [5] Iván Cantador, Ignacio Fernández-Tobías, Shlomo Berkovsky, and Paolo Cremonesi. Cross-domain recommender systems. In *Recommender systems handbook*, pages 919–959. Springer, 2015.
- [6] Sonny Han Seng Chee, Jiawei Han, and Ke Wang. Rectree: An efficient collaborative filtering method. In *International Conference on Data Warehousing and Knowledge Discovery*, pages 141–151. Springer, 2001.
- [7] Tiago Cunha, Carlos Soares, and André CPLF de Carvalho. Selecting collaborative filtering algorithms using metalearning. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pages 393–409. Springer, 2016.
- [8] Antonio Davila, George Foster, and Mahendra Gupta. Venture capital financing and the growth of startup firms. *Journal of business venturing*, 18(6):689–708, 2003.
- [9] Xiaofang Deng, Leilei Wu, Xiaolong Ren, Chunxiao Jia, Yuansheng Zhong, and Linyuan Lü. Inferring users’ preferences through leveraging their social relationships. In *IECON 2017-43rd Annual Conference of the IEEE Industrial Electronics Society*, pages 5830–5836. IEEE, 2017.
- [10] Gary Dushnitsky and Michael J Lenox. When do firms undertake r&d by investing in new ventures? *Strategic Management Journal*, 26(10): 947–965, 2005.
- [11] Martin J Eppler and Jeanne Mengis. The concept of information overload: A review of literature from organization science, accounting, marketing, mis, and related disciplines. *The information society*, 20(5):325–344, 2004.
- [12] Ignacio Fernández-Tobías, Matthias Braunhofer, Mehdi Elahi, Francesco Ricci, and Iván Cantador. Alleviating the new user problem in collaborative filtering by exploiting personality information. *User Modeling and User-Adapted Interaction*, 26(2-3):221–255, 2016.
- [13] Zeno Gantner, Lucas Drumond, Christoph Freudenthaler, Steffen Rendle, and Lars Schmidt-Thieme. Learning attribute-to-feature mappings for cold-start recommendations. In *2010 IEEE International Conference on Data Mining*, pages 176–185. IEEE, 2010.



- [14] Tomer Geva and Jacob Zahavi. Empirical evaluation of an automated intraday stock recommendation system incorporating both market data and textual news. *Decision support systems*, 57:212–223, 2014.
- [15] Tomas Ginevičius, Artūras Kaklauskas, Paulius Kazokaitis, and Jurgita Alchimovienė. Recommender system for real estate management. *Business: Theory and Practice/Verslas: Teorija ir Praktika*, 12(3):258–267, 2011.
- [16] Scott A Golder and Bernardo A Huberman. Usage patterns of collaborative tagging systems. *Journal of Information Science*, 32(2):198–208, 2006.
- [17] Paul Gompers, William Gornall, Steven N Kaplan, and Ilya A Strebulaev. How do venture capitalists make decisions? Technical report, National Bureau of Economic Research, 2016.
- [18] J A Hanley and B J McNeil. The meaning and use of the area under a receiver operating characteristic (roc) curve. *Radiology*, 143(1):29–36, 1982.
- [19] Thomas Hofmann. Latent semantic models for collaborative filtering. *ACM Transactions on Information Systems (TOIS)*, 22(1):89–115, 2004.
- [20] Liangjie Hong, Aziz S Doumith, and Brian D Davison. Co-factorization machines: modeling user interests and predicting individual decisions in twitter. In *Proceedings of the sixth ACM international conference on Web search and data mining*, pages 557–566. ACM, 2013.
- [21] FO Isinkaye, YO Folajimi, and BA Ojokoh. Recommendation systems: Principles, methods and evaluation. *Egyptian Informatics Journal*, 16(3):261–273, 2015.
- [22] Maciej Kula. Metadata embeddings for user and item cold-start recommendations. In Toine Bogers and Marijn Koolen, editors, *Proceedings of the 2nd Workshop on New Trends on Content-Based Recommender Systems co-located with 9th ACM Conference on Recommender Systems (RecSys 2015), Vienna, Austria, September 16-20, 2015.*, volume 1448 of *CEUR Workshop Proceedings*, pages 14–21. CEUR-WS.org, 2015. URL <http://ceur-ws.org/Vol-1448/paper4.pdf>.
- [23] Hao Liao, Manuel Sebastian Mariani, Matúš Medo, Yi-Cheng Zhang, and Ming-Yang Zhou. Ranking in evolving complex networks. *Physics Reports*, 689:1–54, 2017.
- [24] Blerina Lika, Kostas Kolomvatsos, and Stathes Hadjiefthymiades. Facing the cold start problem in recommender systems. *Expert Systems with Applications*, 41(4):2065–2073, 2014.
- [25] Greg Linden, Brent Smith, and Jeremy York. Amazon. com recommendations: Item-to-item collaborative filtering. *IEEE Internet computing*, (1):76–80, 2003.
- [26] Linyuan Lü, Matúš Medo, Chi Ho Yeung, Yi-Cheng Zhang, Zi-Ke Zhang, and Tao Zhou. Recommender systems. *Physics Reports*, 519(1):1–49, 2012.
- [27] Carmen Martinez-Cruz, Carlos Porcel, Juan Bernabé-Moreno, and Enrique Herrera-Viedma. A model to represent users trust in recommender systems using ontologies and fuzzy linguistic modeling. *Information Sciences*, 311:102–118, 2015.
- [28] Colin Mason and Matthew Stark. What do investors look for in a business plan? a comparison of the investment criteria of bankers, venture capitalists and business angels. *International small business journal*, 22(3):227–248, 2004.
- [29] Koji Miyahara and Michael J Pazzani. Collaborative filtering with the simple bayesian classifier. In *Pacific Rim International conference on artificial intelligence*, pages 679–689. Springer, 2000.
- [30] Preeti Paranjape-Voditel and Umesh Deshpande. A stock market portfolio recommender system based on association rule mining. *Applied Soft Computing*, 13(2):1055–1063, 2013.
- [31] Neil Rubens, Mehdi Elahi, Masashi Sugiyama, and Dain Kaplan. Active learning in recommender systems. In *Recommender systems handbook*, pages 809–846. Springer, 2015.
- [32] Aidin Salamzadeh and Hiroko Kawamorita Kesim. Startup companies: life cycle and challenges. In *4th International conference on employment, education and entrepreneurship (EEE), Belgrade, Serbia*, 2015.
- [33] Badrul Munir Sarwar, George Karypis, Joseph A Konstan, John Riedl, et al. Item-based collaborative filtering recommendation algorithms. *Www*, 1:285–295, 2001.
- [34] J Ben Schafer, Joseph A Konstan, and John Riedl. E-commerce recommendation applications. *Data mining and knowledge discovery*, 5(1-2):115–153, 2001.
- [35] Andrew I Schein, Alexandrin Popescul, Lyle H Ungar, and David M Pennock. Methods and metrics for cold-start recommendations. In *Proceedings of the 25th annual international ACM SIGIR conference on Research and development in information retrieval*, pages 253–260. ACM, 2002.
- [36] Mingsheng Shang, Zike Zhang, Tao Zhou, and Yicheng Zhang. Collaborative filtering with diffusion-based similarity on tripartite graphs. *Physica A-statistical Mechanics and Its Applications*, 389(6):1259–1264, 2010.
- [37] Thomas Stone, Weinan Zhang, and Xiaoxue Zhao. An empirical study of top-n recommendation for venture finance. In *Proceedings of the 22nd ACM international conference on Conference on information & knowledge management*, pages 1865–1868. ACM, 2013.
- [38] Zakris Strömquist. Matrix factorization in recommender systems: How sensitive are matrix factorization models to sparsity?, 2018.
- [39] Xiaoyuan Su and Taghi M Khoshgoftaar. A survey of collaborative filtering techniques. *Advances in artificial intelligence*, 2009, 2009.
- [40] Jiliang Tang, Xia Hu, and Huan Liu. Social recommendation: a review. *Social Network Analysis and Mining*, 3(4):1113–1133, 2013.
- [41] Roberto Torres, Sean M McNee, Mara Abel, Joseph A Konstan, and John Riedl. Enhancing digital libraries with techlens+. In *Proceedings of the 4th ACM/IEEE-CS joint conference on Digital libraries*, pages 228–236. ACM, 2004.
- [42] Hastie Trevor, Tibshirani Robert, and Friedman JH. The elements of statistical learning: data mining, inference, and prediction, 2009.
- [43] Xiwang Yang, Yang Guo, Yong Liu, and Harald Steck. A survey of collaborative filtering based social recommender systems. *Computer Communications*, 41:1–10, 2014.
- [44] Fei Yu, An Zeng, Sébastien Gillard, and Matúš Medo. Network-based recommendation algorithms: A review. *Physica A: Statistical Mechanics and its Applications*, 452:192–208, 2016.
- [45] Yi-Cheng Zhang, Marcel Blattner, and Yi-Kuo Yu. Heat conduction process on community networks as a recommendation model. *Physical review letters*, 99(15):154301, 2007.
- [46] Zi-Ke Zhang, Chuang Liu, Yi-Cheng Zhang, and Tao Zhou. Solving the cold-start problem in recommender systems with social tags. *EPL (Europhysics Letters)*, 92(2):28002, 2010.

- [47] Zi-Ke Zhang, Tao Zhou, and Yi-Cheng Zhang. Personalized recommendation via integrated diffusion on user-item-tag tripartite graphs. *Physica A: Statistical Mechanics and its Applications*, 389(1):179–186, 2010.
- [48] Xiaoxue Zhao, Weinan Zhang, and Jun Wang. Risk-hedged venture capital investment recommendation. In *Proceedings of the 9th ACM Conference on Recommender Systems*, pages 75–82. ACM, 2015.
- [49] Tao Zhou, Jie Ren, Matúš Medo, and Yi-Cheng Zhang. Bipartite network projection and personal recommendation. *Physical Review E*, 76(4):046115, 2007.
- [50] Dávid Zibriczky. Recommender systems meet finance: A literature review. In *International Workshop on Personalization and Recommender Systems in Financial Services*, 2016.

## Appendix A. Different training/testing set splitting

We test the 85%/15% and 95%/5% splitting of training and testing set to examine the stability of the results. Figure A.8 shows the similar relative performance, which is in agreement with Figure 6 (c). Figure A.9 is in agreement with Figure 7 (a).

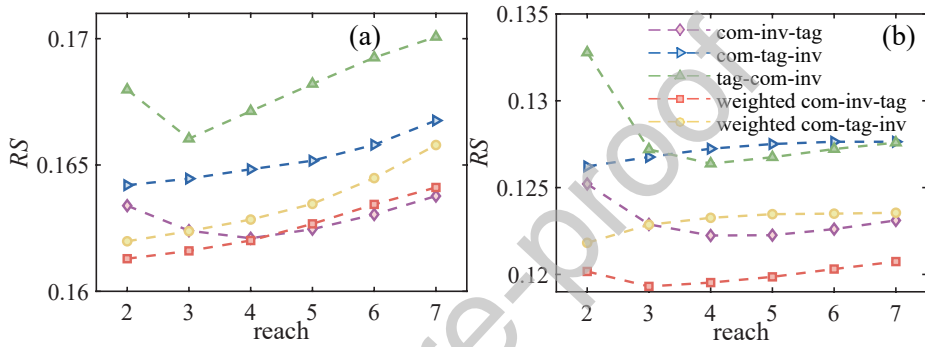


Figure A.8: Performance of the Probs algorithm applied to five different tripartite network representations, for different sizes of the testing set. In (a), the proportion of links in the testing set is 15%. In (b), the proportion of links in the testing set is 5%. The methods show a similar relative performance as in Figure 6 (c)

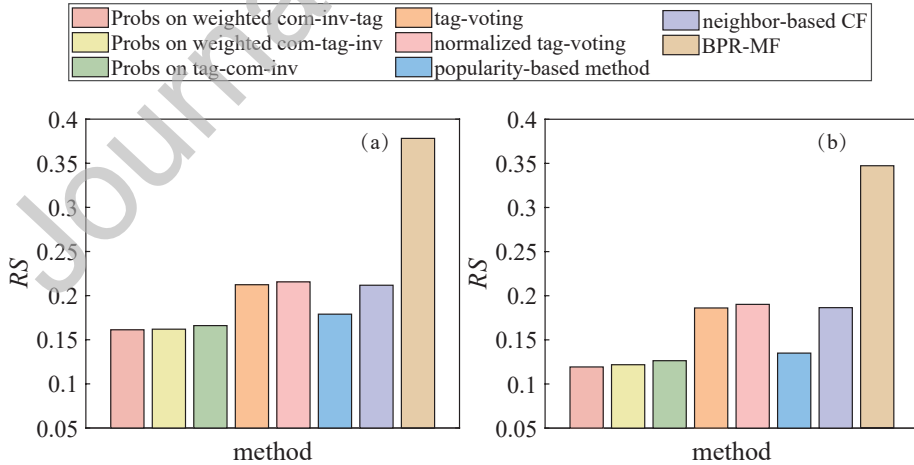


Figure A.9: A comparison of considered recommendation methods with different size of the testing set. In (a), the proportion of links in the testing set is 15%. In (b), the proportion of links in the testing set is 5%. The methods show a similar relative performance as in Figure 7 (a).

21 November 2019

Dear Editor,

We certify that there is no conflict of interest to report.

Yours sincerely,

Shuqi Xu    Qianming Zhang    Linyuan Lü    Manuel Sebastian Mariani

## **Author contributions**

### **Author 1: Shuqi Xu**

Data Curation, Investigation, Software, Formal analysis, Writing - Original Draft, Writing - Review & Editing

### **Author 2: Qianming Zhang**

Conceptualization, Methodology, Data Curation, Validation, Writing - Original Draft

### **Author 3: Linyuan Lü**

Conceptualization, Methodology, Supervision, Writing - Review & Editing

### **Author 4: Manuel Sebastian Mariani**

Writing - Original Draft, Writing - Review & Editing